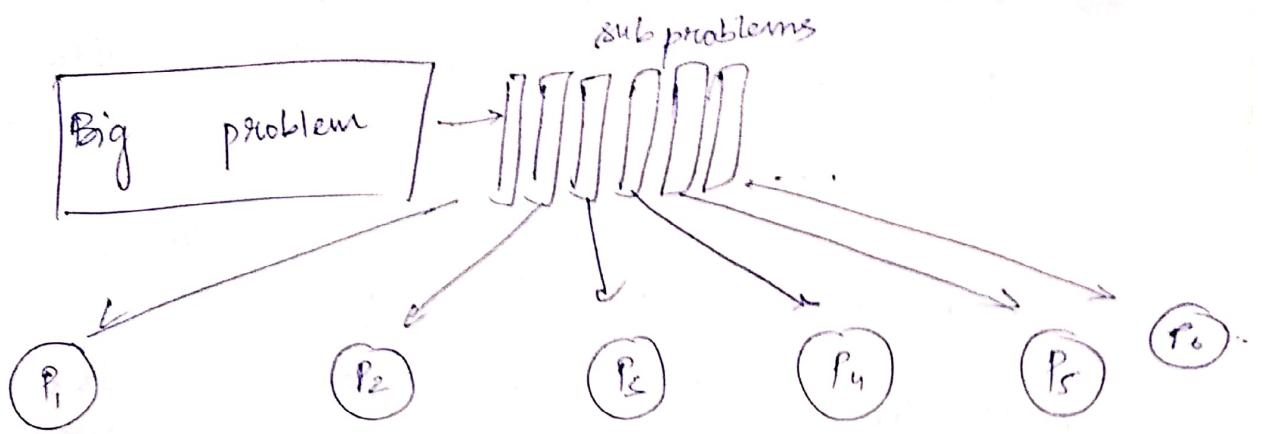
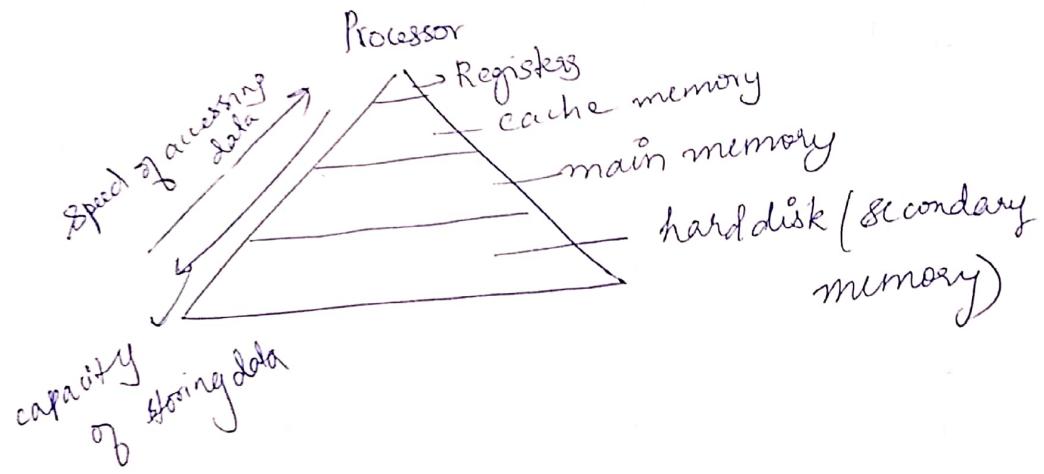


## Parallel computing



## memory hierarchy

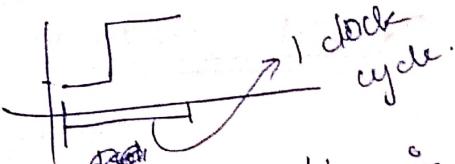


$$\text{Performance}_x = \frac{1}{\text{execution Time}_x}$$

$x$  is  $n$  times faster than  $y$

$$\frac{\text{Performance}_x}{\text{Performance}_y} = \frac{\text{Execution}_Y}{\text{Execution}_x} = n$$

Clock cycles



→ It's better to express execution time in ~~seconds~~ cycles rather than in seconds

→ clock ticks → start and end of a cycle.

→ cycle time = time between ticks = seconds per cycle

→ ~~cycle ticks~~ clock rate = cycles per second.

→ 1 Hz = 1 cycle/sec, 1 MHz =  $10^6$  cycles/sec

Performance Eq(1)

$$\text{Clock cycle time} = \frac{1}{\text{clock rate}}$$

$$\frac{\text{seconds}}{\text{program}} = \frac{\text{cycles}}{\text{program}} \times \frac{\text{seconds}}{\text{cycle}}$$

cpu execution time = (cpu clock cycles for a program)  $\times$  clock cycle time for a program

Performance :- The capacity to reduce the execution time of a program with some additional resources.

computer performance is based upon.

(i) Response Time :- the time it takes for the program to start and finish, it includes accessing disk and memory, waiting for I/O and other processes, OS overhead.

(ii) Throughput :- It talks about how much amount of the data is being processed in a given time

or

It talks about how ~~many~~ many processes or jobs can the machine run at once.

Q) Execution time initial =  $10^8$

Clock rate =  $4 \times 10^9$  cycles/sec

Execution time new =  $6 \text{ s}$

$1 \text{ ns} = 10^{-9} \text{ seconds}$

$\therefore 10^8 = \frac{\text{CPU clock cycles A}}{4 \times 10^9 \text{ cycles/sec}}$

$\text{CPU} = 1 \text{ ns}$

-----  
CPU clock cycles A =  $4 \times 10^9 \times 10^8$  cycles

CPU clock cycles B =  $1.2 \times \text{CPU clock cycles A}$

$\therefore \text{Execution time} = 6 = \frac{\text{CPU clock cycles B}}{\text{Clock rate}}$

Clock rate =  $\frac{2 \times 4 \times 10^9 \times 10^8}{6}$

$= 8 \times 10^9$

$= 8 \text{ GHz}$

## Performance Eq(2)

CPU execution time

$$\text{units}^{-1} = \frac{\text{seconds}}{\text{program}}$$

Instruction count for a program

CPU clock cycles

$\times$  average CPI

clock cycles time

$$\therefore \boxed{\text{CPU } E_x = IC \times \text{CPI} \times \text{clock cycles}}$$

CPI  $\rightarrow$  clock cycles per instruction

$$\text{CPU execution time} = \frac{\text{seconds}}{\text{program}} = \frac{\text{Instructions}}{\text{program}} \times \frac{\text{cycles}}{\text{instruction}} \times \frac{\text{seconds}}{\text{cycle}}$$

Q)

$$\frac{\text{CPU execution A}}{\text{CPU execution B}} = \frac{7 \times 250 \times 20}{7 \times 500 \times 1.2}$$

$$\Rightarrow \frac{\text{CPU execution B}}{\text{CPU execution A}} = \frac{1.2}{7} = \frac{\text{Performance B}^A}{\text{Performance A}^B}$$

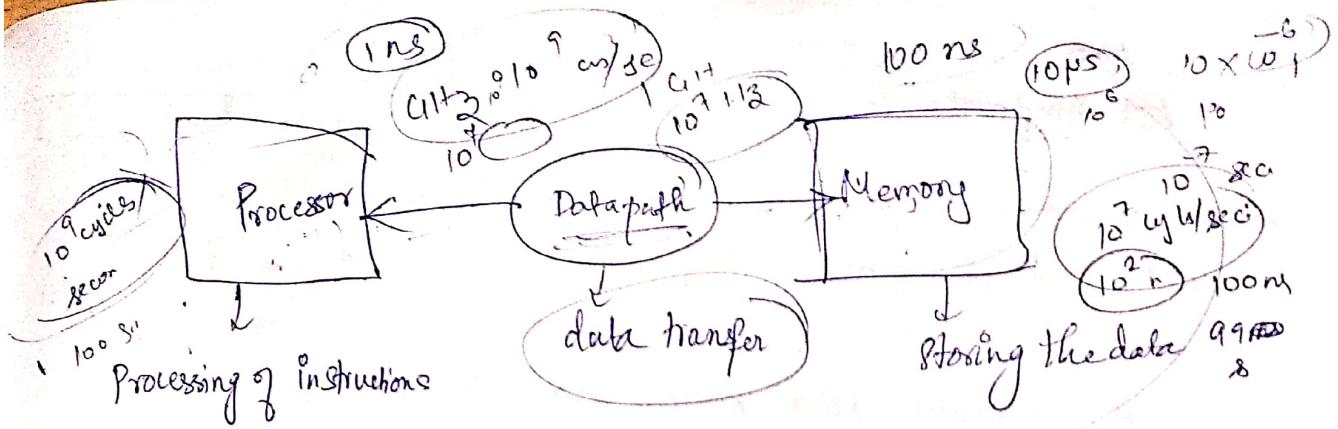
$$\boxed{\text{CPU clock cycles} = \frac{\text{Instruction count for a Program}}{\text{average CPI}}} \times \text{clock cycles}$$

$$10 = (5) \times c \\ = \text{CPI}_1 = 2$$

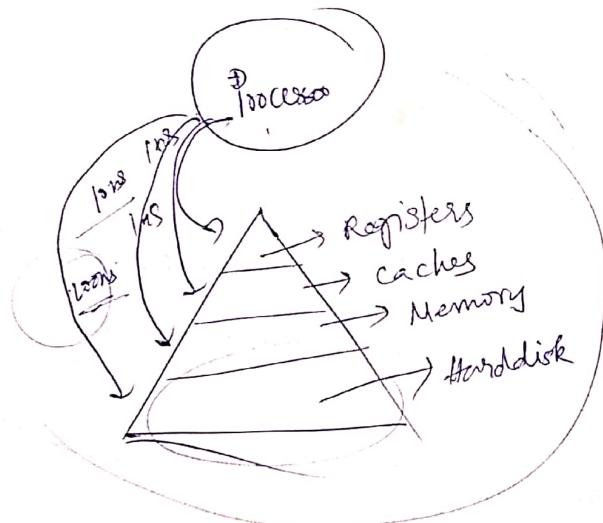
$$9 = (6) \times \text{CPI}_2 \\ \text{CPI}_2 = 1.5$$

CPU execution time :- It refers to the time in which CPU processes the task or executes the given task (excludes I/O, running other programs)

this is simply referred to as CPU time



- Memory is very slow compared to processor
- Bottlenecks occur as the rate at which memory sends data is less than the rate at which the processor processes data.



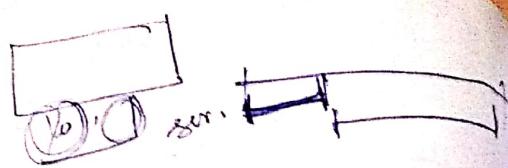
Multiplicity → more no. of processors, more, data paths, more memory.

Implicitly → compilers will take care of parallelism

Indirectly → parallelism is taken by the programmer.

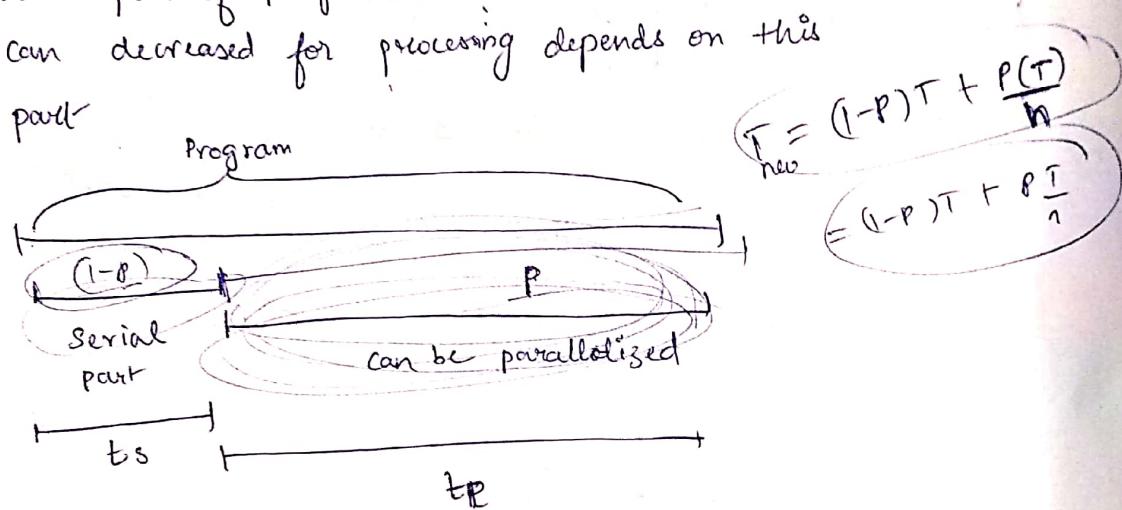
~~Area~~

## And Amdahl's law



→ In a ~~program~~ program there will be serial processes and some part of the serial process which can be parallelized.

→ So the Amdahl's law states that there will be a performance improvement when we parallelize the part of program, and the time that can decrease for processing depends on this part.



Speedup =  $\frac{\text{Performance for the task using the enhancement}}{\text{Performance for the task without using enhancement}}$  *when possible*

(or)

=  $\frac{T}{(1-P)T + PT}$

= Execution time taken for a task without using enhancement

Execution time taken for a task with using enhancement when possible

Let the p be the fraction of the task to be enhanced.

$$T_{\text{old}} = T$$

$$(T_{\text{new}} = T(1-P) + \frac{TP}{n})$$

$$T_{\text{new}} = T_{\text{old}}(1-P) + \frac{T_{\text{old}}P}{n}$$

$$\frac{T_{old}}{T_{new}} = \frac{1}{(1-p) + \frac{p}{n}}$$

→ here  $n$  is the no. of processors  
 →  $p$  is fraction of the task that can be parallelized

$$\boxed{\text{SpeedUp} = \frac{1}{(1-p) + \frac{p}{n}}}$$

$$S_U = \frac{T_{old}}{T_{new}}$$

$T_{old} = 100$

$$20/100 = \cancel{100} + \frac{80}{n}$$

$$S_U = \frac{T_{old}}{T_{new}} \quad T_{new} = 20.8$$

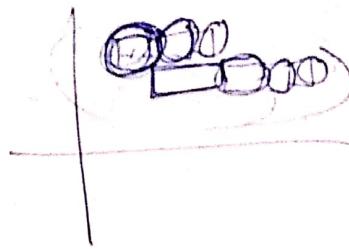
$$ex \rightarrow \frac{\text{Speedup}}{T_{old}} = 10 \quad n = 5$$

$$T_{new} = 5 + \frac{5}{5} = 6$$

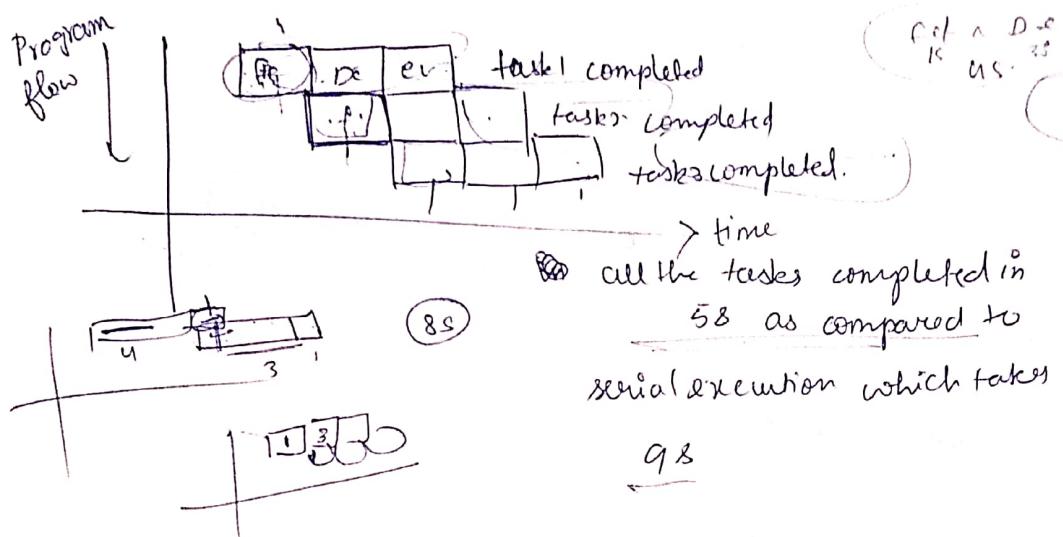
$$\frac{1}{5+5} = \frac{1}{6}$$

$$\therefore \text{Speedup} = \frac{10}{6} = 1.6$$

- Q) How processor relies on different process to achieve parallelism and achieve better performance.
- Pipelining
  - Superscaling
  - Very long instruction word



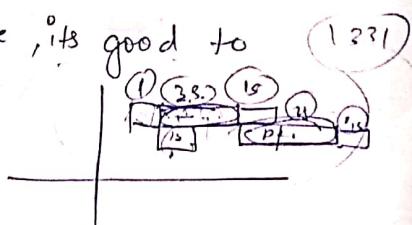
Pipelining :- It overlaps various stages of instruction during execution to achieve performance.



→ Processors have relied on pipelining technique to improve execution rate.

→ To increase speed of a single pipeline, one would break down the tasks into smaller units.

→ To further improve the execution rate, it's good to use multiple pipelines.



### Problems in Pipelines

→ The pipeline gets slowed down due to slowest step or stage.

• To solve this, a conventional processor relies on a very deep pipeline (upto 20 stages).

→ Even then as every instruction not only consists of add or mult, it may even have jump conditions based on the if and else constructs, this jump requires an accurate branch prediction.

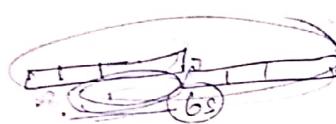
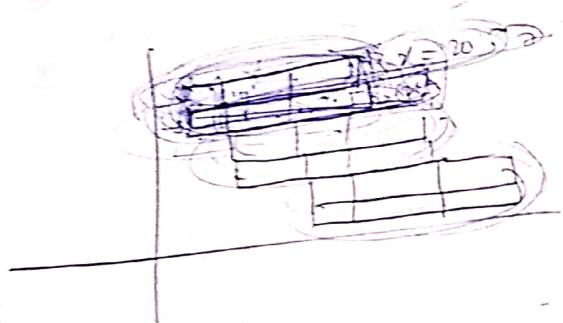
- If there is a misprediction, then the penalty of such prediction grows through the depth of the pipeline and a large number of instructions have to be flushed.

### Superscalar

- one simple way of dealing with problems in pipelining is to

use multiple pipelines

- same time 2-2 instructions will be fetched (e.g.) in pipelining method
- same time single instruction being fetched, but here 2 instructions are fetched at a time.



$$\begin{aligned} n &= 2 \rightarrow \\ m_d &= 2 \times 1 \rightarrow 3 \\ n &= 3 \rightarrow \\ m_d &= 3 \times 2 \rightarrow 6 \\ y &= 2k + l - 1 \\ y &= n_1 - 1 \end{aligned}$$

### Issues in Superscalar

- not all functions are kept busy at all times.
- ~~idle resources~~ if during a cycle no functional units are utilized then it refers to the vertical waste
- if during a cycle, only some of the functional units are utilized, this refers to as a horizontal waste.
- data dependencies may waste some functional units.

Scheduling of instructions is determined by a number of factors

- true data dependency - the result of one op is an ifp to

other ops

- Resource dependency - one resource required by 2 op

- Branch dependency - for conditional branch instruction destination is only known at the point of execution

- In simpler model, instructions can be issued only in the order in which they are encountered, i.e., if the second instruction has a data dependency with the first, then the second instruction is not issued along with the first instruction. This is called in-order issue.
- In a more aggressive model, if the ~~1<sup>st</sup>~~ and 2<sup>nd</sup> instruction has a data dependency with 1<sup>st</sup> and there is a 3<sup>rd</sup> instruction which doesn't have any dependency with 1<sup>st</sup> and 2<sup>nd</sup> instruction, then in the first clock cycle the 1<sup>st</sup> and 3<sup>rd</sup> instruction can be clubbed, this is called a dynamic issue.
- There is a hardware tool called scheduler which looks into the long queue and selects appropriate instructions to execute <sup>concurrently</sup>. To overcome the dependency issues present in the based on the superscalar we use VLIW <sup>above dependencies</sup>

### VLIW (Very long instruction word)

- These processors came into existence ~~as~~ in order to address the complexity of superscalar processor and its scheduler.
- VLIW processors rely on compile time analysis to analyse what instructions to bundle together that can be executed concurrently.
- These instructions are ~~not~~ packed and dispatched together as a single instruction long word hence VLIW

## Issues with ~~the~~ VLIW

- compiler has a bigger context from which it has to analyse
- compilers however don't have runtime info
- branch and memory prediction is more difficult
- VLIW performance is highly dependent on the processor
- A typical VLIW implements 4-way or 8-way parallelism

## Limitations of Memory system Performance

- memory speed is less as compared to processor speed, hence bottlenecks occur
- memory system performance is largely captured by two parameters, latency and bandwidth
- latency is the time taken by ~~the memory~~ from the issuing of a memory request to the time the data is available at processor.

for example, to fetch a ~~value~~ value of A, the processor sent a request at the 0<sup>th</sup> second and the data of A is available at processor at 2nd second  
∴ the latency here is 2s.

bandwidth → It's the rate at which the data is being pumped to the processor.

$$\text{Ex:- processor speed} = 1 \text{ GHz} (1 \text{ ns clock})$$

$$\text{CPU clock} = \frac{1}{\text{clock time}}$$

DRAM latency = 100ns (no caches)

$$10^{-9} \text{ s}$$

$$\text{peak performance} = (\text{GFlops})$$

$$\text{CPU speed in GHz} \times \left( \frac{\text{CPU instruction}}{\text{per cycle}} \right)$$

$$\text{GFlops} = \frac{\text{cycles}}{\text{second}} \times \frac{\text{instructions}}{\text{cycle}} = \frac{\text{instructions}}{\text{second}}$$

processor executes 4 instructions per cycle.

$$\therefore \text{Peak performance} = 4 \text{ GFlops} = 4 \times 1 \text{ GHz} \\ = 4 \text{ GFlops}$$

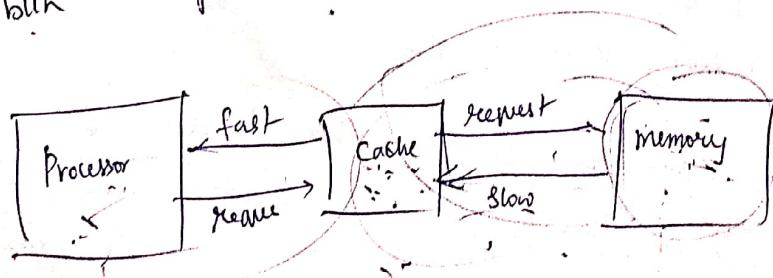
But due to the memory latency, the processor has to wait for  $\frac{100}{\text{clock}}$  cycles to process the data.

$\therefore$  bandwidth or ~~bandwidth~~ It depends on how much fast the memory pumps in the data to the processor, hence

$$\text{bandwidth} = \frac{1}{100 \text{ ns}} = \frac{1}{100 \times 10^{-9}} = 10^7 \text{ GFlops} \\ = 10 \times 10^6 \text{ FLOPs} \\ = 10 \text{ MFLOPs}$$

Improving effective memory latency using Caches

→ Caches are small and fast memory elements built between the processor and the DRAM



→ If a piece of data is repeatedly used, the effective latency of the memory system can be reduced by cache.

- The fraction of data references satisfied by the cache is called cache hit ratio.
- Initially some things are not present on the cache, in the initial request the cache will not have the data and hence it's a miss, therefore cache will get the data from main memory and store it, in the next ~~request~~ request the cache will be able to provide the data to processor and hence it's a hit.

$$\text{hit ratio} = \frac{\text{no. of hits}}{\text{no. of hits} + \text{no. of misses}}$$

### EXAMPLE & WORKS

consider the bifer page example and there is a cache of 32 KB

so first we need to fetch the ~~two~~ two matrices from the main memory. so for that there will be some time it takes

lets calculate it (delay)  $\rightarrow 2 \times 32 \times 32 \times 100 \text{ ns delay}$

matrices matrix

$$\Rightarrow 2 \times 32 \times 32 \times 100 \times 10^{-9}$$

$$640000$$

$$10^9$$

$$\approx 200 \mu\text{s}$$

$$200$$

$$\frac{1}{640000}$$

$$\begin{aligned} & \cancel{X = S \times D} \\ & \cancel{\text{rate} = X / \text{time}} \end{aligned}$$

→ now multiplying 2 matrices will take  $2n^3$  operations.

∴ the total operation  $\rightarrow 2 \times 32^3 \approx 64 \text{ K operations}$

∴ the processor performs  $4 \text{ instruction/operations per cycle}$

$$\therefore \text{total no. of cycle required by the processor} = \frac{64 \text{ K}}{4} = 16 \text{ K cycles}$$

(for 16 ps)

$$= 16 \text{ K cycles} / \text{s.}$$

now the peak computation rate =  $\frac{\text{CPU speed} \times \text{no. of operations}}{\text{instructions}}$   
 or operations per cycle

$$\text{total delay} = 200\ \mu\text{s} + 16\ \mu\text{s}$$

$$= 216\ \mu\text{s}$$

$$\text{CPU speed} = \frac{1}{216\ \mu\text{s}}$$

∴ peak computation rate =  $\frac{64 \times 10^9 \text{ instructions}}{64 \times 10^9 \text{ instruction/cycle}} \times \frac{1}{216\ \mu\text{s}}$

$$= \frac{64 \times 10^9}{216 \times 10^6}$$

$$= \frac{16 \times 10^3}{216}$$

$$= 0.2962 \times 10^9$$

$$= 296 \times 10^6 \text{ F}$$

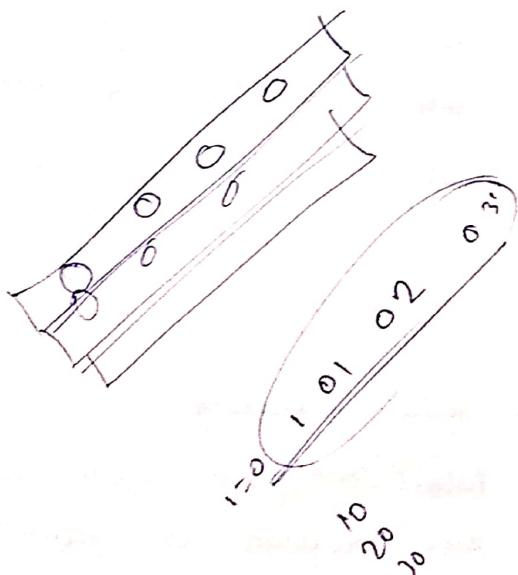
$$\Rightarrow 296 \text{ MFlops}$$

### Impact of memory bandwidth

- memory bandwidth is determined by the bandwidth of the memory bus as well as the memory units.
- memory can be improved by increasing the size of memory blocks.
- spatial locality → storage of consecutive words in the cache
- temporal locality → storage of nonconsecutive words in the cache.

## Some other ways of reducing memory latency

- prefetching → It's a technique used by computer processors to boost execution performance by fetching instructions or data from their original storage in slower memory to a faster local memory even before they are needed.
- multi-threading → creating different threads ~~at the same time~~ in order to execute different operations which are not dependent on each other.
- spatial locality → It's a technique in which we bring in all the data into the cache ~~from the main memory~~ from the main memory.



## Explicitly parallel programs

→ this type of parallelism is controlled by the programmers.

→ a parallel program must exhibit concurrency and interaction b/w the concurrent subtasks

control structure

communication model

→ parallelism can be expressed at various levels of granularity.

from instruction level to processes level

### Control structure of parallel programs

→ processing units parallel computers either operate under a ~~single~~ centralised control (or) → work independently

→ If there is a single control unit, that sends ~~instructions~~ instructions to various processors that work on different data, this model is referred to as SIMD → Single instruction <sup>stream</sup>, multiple data stream.

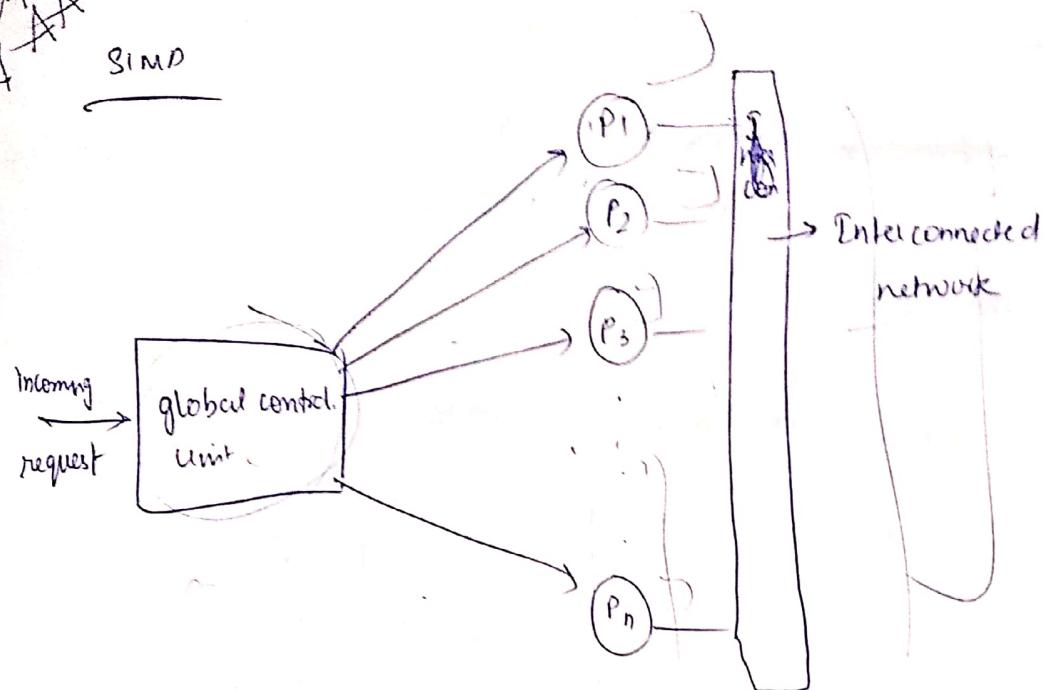
→ If ~~multiple~~ <sup>multiple</sup>

MIMD → each processor has its own control unit, each processor can execute different instructions on different data items

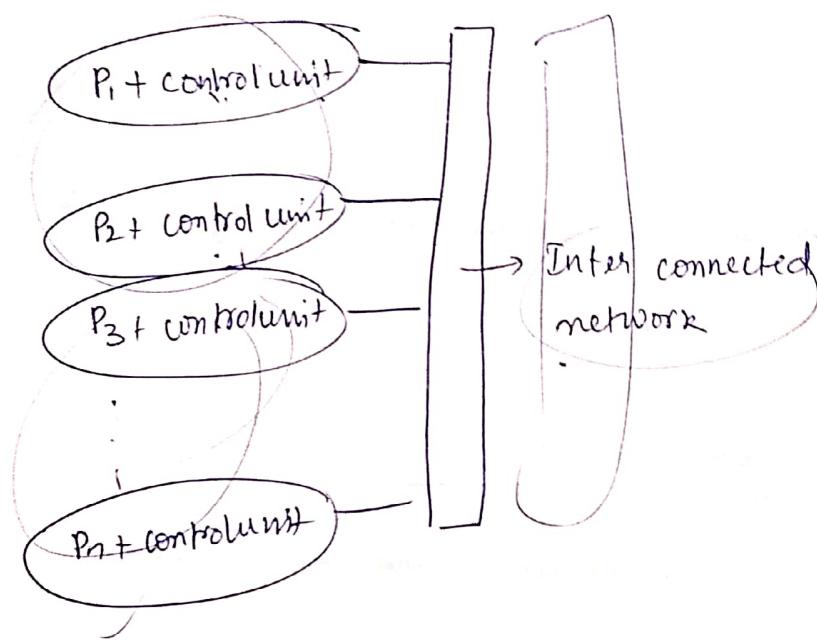
multiple instruction stream, multiple data stream,

A A

SIMD



MIMD



SIMD

- example MAC IV, MPP, DAP
- sometimes we ~~want~~ don't want some operations on certain data items, an activity mask determines whether a processor should participate in the computation or not.
- single program, on different processors

## MIMD

- multiple programs on multiple processors.
- single program multiple data streams.
- Sun ultra servers, IBM servers

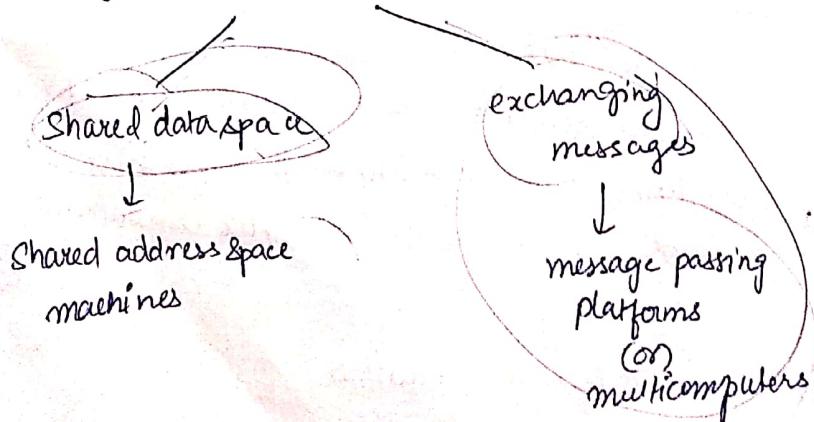
## SIMD vs MIMD

- SIMD requires less hardware than MIMD because executing on the same instruction and also there is one global control unit whereas on MIMD each processor has its own control unit hence the hardware complexity increases.
- SIMD are specially designed hence they tend to be expensive as compared to MIMD
- Not all the applications are naturally suitable with SIMD

## Communication Model

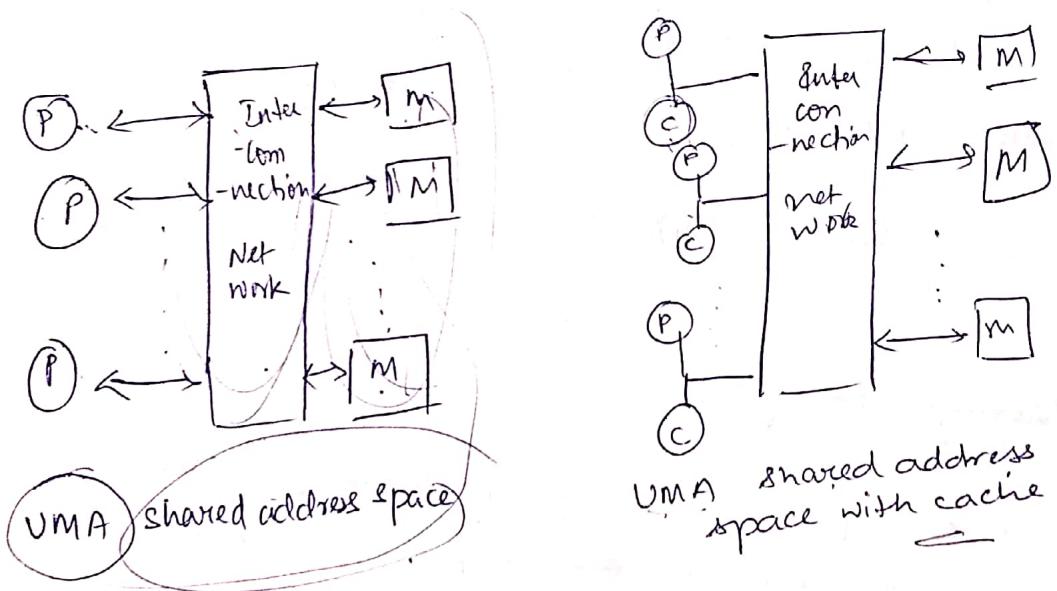
How does the communication b/w the processors.

- two forms of data exchange

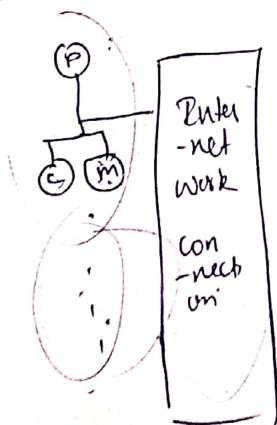


## Shared address space platforms

- partial or all the data is ~~is~~ available to all processors
- processors interact by modifying the data in shared address space
- If the time taken by the processor to access any memory in the global ~~is~~ space or local, then it's referred uniform memory access (UMA) and else it's referred as ~~is~~ uniform ~~access~~ memory access (NUMA)



UMA shared address space with cache



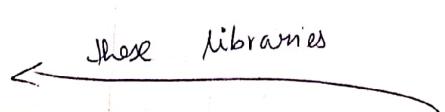
2 NUMA shared-address space with local memory

## NUMA and UMA

~~UMA~~

- The distinction b/w NUMA and UMA is very much important as in NUMA each and every processor has its own cache and ~~no~~ local memory whereas on the other side in UMA, ~~less or shared~~ the processors have shared memory, hence an algorithm should be defined in such a way that the processor can get the data efficiently
- programming of these platforms is easier since read and write are visible to other processor.
- There should be a co-ordination of read-write data to the shared address space as there might be conflicts when one processor is reading from the memory and at the same time another process might be writing into it.
- Cache coherence problems

## Message Passing platforms

- We are passing messages from one processor to another as one processor's output may serve as an input to other one.
- These platforms are especially used in cluster work stations.
- Ex MPI, PVM
- Programming of these platforms are done through  these libraries

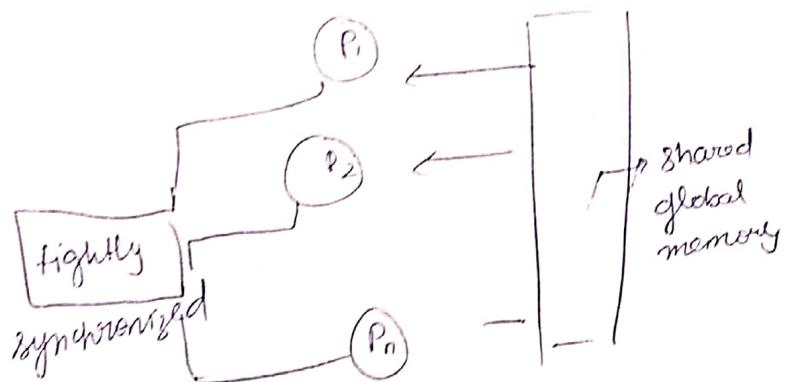
## MP & V Shared-Address Space

- MP requires little hardware support, other than network whereas shared-address-space requires a lot of hardware support (i.e) there should be controllers which make sure that there are no conflicts while reading and writing data to it.

## Physical Organization of Parallel platforms

PRAM → Parallel random access machine

- PRAM consists of  $p$  processors and they all have ~~all~~ global access to the global shared memory.
- PRAM → all the processors have a common clock but can execute different instructions in each cycle



Depending on how simultaneously it access the memory PRAMs are classified as

- exclusive read - exclusive write ERW PRAM
- concurrent read - exclusive II CREW PRAM
- Exclusive read - concurrent write CRCW PRAM
- concurrent II - concurrent II CRCW PRAM

## Concurrent



Common values are allowed to write



Arbitrary process will be given ability to write



Priority is given to processes to write



Sum of all the data is written <sup>in</sup> to the memory

- processors and memories are connected via switches.
- these switches must operate in  $O(1)$
- interconnection networks carry data <sup>between</sup> processors and to memory
- interconnects are made of links and switches
- static interconnects are interconnect which connect point to point there are also direct connections
- dynamic ~~as~~ network are built using switches and links. these are indirect networks

## Network Topologies

### → Buses

- distance b/w any two nodes is  $O(1)$
- all entities involved share the same bus

### Crossbars

- It's a pxm network

## Thread

- A thread is a single sequential flow of control within a program.
- multiple threads can be executed in parallel on computer systems.
- It occurs by time slicing, ~~as~~ time-slicing is based on non-priority scheduling.
- threads typically share memory and other resources directly.
- The communication overhead is less as they use the shared memory space.
- user level → user level will take care.
- kernel level → operating system will take care
- the threads can be used via Pthread header file and can be compiled using gcc itself.
- Stack will be created for each thread
- threads will share the shared-address-space and have their own memory, hence it leads to poor performance.
- Threads are <sup>software</sup>, portable because every processor are capable of ~~doing~~ handling the function of threads.
- Inherent support for latency hiding → one thread is independent of others, for example while one thread is doing I/O operations, the other thread can do some other operation.

→ scheduling and load balancing

→ ease of programming

### Thread creation and termination

#### Thread safety

→ it's the ability to execute multiple threads simultaneously without clobbering thread data or creating race conditions.

#### Why pthreads?

→ performance gains

→ cost of managing a process is more than managing a thread as it requires much fewer system resources.

→ on modern, CPU's, pthreads are suitable for parallel programming

#### pthread API

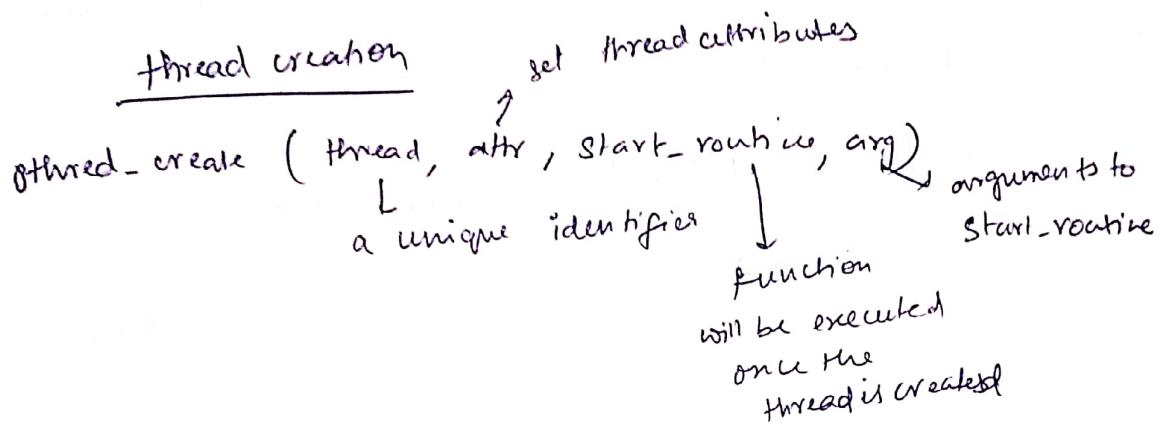
1) thread management :- functions that directly work on threads

- creating
- detaching
- joining

2) mutex :- functions that deal with synchronization

called a mutex, this includes creating, destroying, locking, unlocking mutexes

- ↳ conditional variables → functions that address communications b/w threads that share a mutex.
- this class includes functions to create destroy wait and signal based on ~~mutexes~~ specified variables



### thread termination

- the thread is terminated when it's returned to main from the start routine
- via pthread-exit() → this is explicit
- one thread is cancelled by other thread pthread-cancel()
- the entire main process is exited ~~due to~~ due to exit routine

pthread-exit is used explicitly exit a thread; it's generally called after the thread has completed its work and no longer required to exist

- pthread-join() tells to execute the thread.

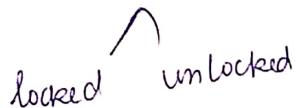
## Synchronization primitives in threads

→ If there is no sync b/w threads, it may result in to unwanted solutions, hence we need synchronization b/w threads.

### Mutual exclusion

```
if (mycost < best cost) } critical sections  
{ best cost = mycost }  
y
```

→ critical segments in pthreads are implemented using mutex locks



→ at any point of time 1 thread can have a mutex lock.

→ ~~a thread~~.

pthread - mutex - lock

pthread - mutex - unlock .

## Producers Consumers using locks

- manage shared memory access
- checking if buffer is full (producers)
- checking if buffer is empty (consumers)

th Num\\_threads = 2

int buffer [10];

int count = 0;

int main () {

    Pthread\_t th[~~Thread~~ Num\_threads]

    for (int i=0; i< num\_threads; i++)

        if (i == 0)

            Pthread\_create(&th[i], NULL, &producer, NULL)

        else

            Pthread\_create(&th[i], NULL, &consumer, NULL)

}

..

void \* producer (void \* args)

{

    int x = rand() % 100;

    buffer [count] = x;

    count++;

}

Types of mutex

→ Normal mutex :- app this mutex deadlocks if a

thread that already has a lock

tries a second lock on it

→ Recursive → single thread is allowed to lock a mutex as many time it wants.

Error check mutex → reports an error when a thread with a lock tries to lock it again.

$$\text{Proc} = 1 \text{ GHz}$$

$$\left( \frac{1}{10^9} \right) 8$$

$$1 \text{ cycle} = \frac{1}{10^9 \text{ s}}$$

$$1 \text{ cy} \rightarrow \left( \frac{1}{10^9} \right) 8$$

$$C = \frac{1}{10^9 \text{ s}}$$

$$100 \rightarrow \frac{1}{10^9}$$



$$n = \frac{100}{10^9}$$

$$n = \frac{1}{10^7 \text{ s}}$$

Cycle / sec

$$8 \times 4 = 8 \text{ words}$$

~~$$8 \times 10^9 \text{ bits}$$~~

$$1 \text{ cy} \left( \frac{1}{10^9 \text{ s}} \right)$$

$$\frac{1}{10^7} 8$$

$10^{-7}$  sec

inv / sec

~~$10^{-7}$~~   $10^7$

$$CP = \frac{\text{cycle}}{\text{second}}$$

~~number of operations~~